Evaluating the Robustness of Watermark-based Detection of Algenerated Content

Neil Gong

Department of Electrical and Computer Engineering Department of Computer Science (secondary appointment) Duke University 10/16/2023

Ethical Concerns of Al-generated Content

- Harmful content
- Disinformation and propaganda campaigns
- Teaching and education

Content Moderation for Generative AI

- Preventing generation of harmful content: safety filters
 - Deployed by many GenAl services
 - Yang et al. "SneakyPrompt: Evaluating Robustness of Text-to-image Generative Models' Safety Filters". *Arxiv*, 2023.

- Detecting Al-generated content
 - Watermark-based detection of AI-generated images is deployed by Google, Stability AI, OpenAI, etc.
 - Jiang et al. "Evading Watermark based Detection of AI-Generated Content". In ACM Conference on Computer and Communications Security (CCS), 2023.

Detecting Al-generated Content

Passive detection

- Key idea: leverage artifacts in AI-generated content
- High false positives/negatives

- Watermark-based detection
 - Multiple companies have deployed such detector
 - This talk
 - Al-generated images

Image Watermarking

- Three components
 - Watermark (bitstring)
 - Encoder
 - Decoder



Watermarking Methods

- Non-learning-based
 - Encoder and decoder are handcrafted based on heuristics
 - Not robust to common post-processing
 - JPEG compression, Gaussian noise, Gaussian blur, Brightness/Contrast



Watermark used by Stable Diffusion

- Learning-based
 - Encoder and decoder are neural networks
 - Believed to be robust due to adversarial training

Standard vs. Adversarial Training



Standard training

Standard vs. Adversarial Training



Adversarial training

Watermark-based Detection



How to Set Detection Threshold τ ?

Achieve a desired False Positive Rate (FPR)



Double-tail Detector





Detector Deployment Scenarios

- Detection-as-a-service
 - Provider of GenAI service also provides detection service
- End-user detection
 - Detector as end-user app
 - Mobile app
 - Browser plugin
- Public detection
 - Publicly release decoder and watermark
 - Individuals can personalize au depending on desired FPR
- Third-party detection
 - GenAI provider shares decoder and watermark with selected third parties
 - E.g., Google \rightarrow Twitter

Threat Model

- White-box setting
 - Attacker has access to decoder
 - Aim to evade detector with any τ > 0.5
- Black-box setting
 - Attacker has access to detector API
 - Aim to evade a specific detector with an unknown au
- Focus on watermark removal
 - Watermark forging/spoofing is technically the same

One Visualization Example



Adversarial training

White-box Setting

- Given a watermarked image I_w
- Add minimal perturbation δ
- s.t. each bit of decoded watermark flips

$$\min_{\delta} ||\delta||_{\infty}$$
s.t. $D(I_{w} + \delta) = \neg D(I_{w})$
Decoder

Decouer

Flip each bit

Guaranteed to evade single-tail detector

Can be detected by double-tail detector

White-box Setting

- Intuition: non-watermarked images have bitwise accuracy pprox 0.5
- Add minimal perturbation to make bitwise accuracy ≈ 0.5
 - Perturbed watermarked image indistinguishable with non-watermarked ones

$$\begin{split} \min_{\delta} ||\delta||_{\infty} & {}_{\text{Ground-truth watermark}} \\ s.t. & |BA(D(I_{w}+\delta),w)-0.5| \leq \epsilon \end{split}$$

White-box Setting

- Intuition: non-watermarked images have bitwise accuracy pprox 0.5
- Add minimal perturbation to make bitwise accuracy ≈ 0.5
 - Perturbed watermarked image indistinguishable with non-watermarked ones



Theoretical Evaluation

Evasion rate: probability that a perturbed watermarked image is detected as non-AI-generated

Lower bound for single-tail detector: $Pr(m \le \lfloor (\tau - \varepsilon)n \rfloor)$ $m \sim Binomial(n, 0.5)$

Lower bound for double-tail detector: $2 \Pr(m \le \lfloor (\tau - \varepsilon)n \rfloor) - 1$

Empirical Evaluation Results



Our Adversarial Post-processing Adds Smaller Perturbations than Existing Ones



COCO dataset

Adversarial Training Improves Robustness



Take-away Messages

- Learning-based-watermark based detection has good robustness to common post-processing in *non-adversarial settings*
- Broken in the white-box setting in *adversarial settings*
- Adversarial training improves robustness but still insufficient

Black-box Setting



Theoretical Evaluation

Evasion rate is guaranteed to be 1

Empirical Evaluation Results



COCO dataset

Adversarial Training Improves Robustness



Summary and Discussion

- Don't publicly release decoder
 - No white-box attack
- Adversarial training can improve robustness
- But probably still insufficient
 - Dozens-hundreds of queries to evade a black-box detector
 - While maintaining image quality
- Stronger adversarial training?

Acknowledgements

- Restricted access to detector API?
 - Attacker cannot access detector API
 - Transfer attack

Zhengyuan Jiang Jinghuai Zhang